

AD-754 751

MULTIWAY CONTINGENCY TABLE ANALYSIS
APPLIED TO THE CLASSIFICATION OF MULTI-
VARIATE DICHOTOMOUS POPULATIONS

S. Kullback

George Washington University

Prepared for:

Office of Naval Research

9 January 1973

DISTRIBUTED BY:

NTIS

National Technical Information Service
U. S. DEPARTMENT OF COMMERCE
5285 Port Royal Road, Springfield Va. 22151

AD-754751

**MULTIWAY CONTINGENCY TABLE ANALYSIS APPLIED TO THE
CLASSIFICATION OF MULTIVARIATE DICHOTOMOUS POPULATIONS**

by

S. KULLBACK

TECHNICAL REPORT NO. 4

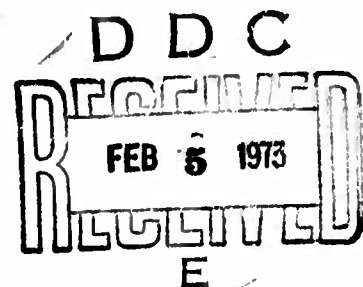
January 9, 1973

PREPARED UNDER CONTRACT N00014-67-A-0214-0015

(NR-042-267)

OFFICE OF NAVAL RESEARCH

Herbert Solomon, Project Director



**Reproduction in Whole or in Part is Permitted for
any Purpose of the United States Government**

**Reproduced by
NATIONAL TECHNICAL
INFORMATION SERVICE
U S Department of Commerce
Springfield VA 22151**

**Approved
for publication and sale; the
distribution is unlimited.**

**DEPARTMENT OF STATISTICS
THE GEORGE WASHINGTON UNIVERSITY
WASHINGTON, D. C. 20006**

Unclassified

Security Classification

DOCUMENT CONTROL DATA - R&D		
(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)		
1. ORIGINATING ACTIVITY (Corporate author) THE GEORGE WASHINGTON UNIVERSITY DEPT. OF STATISTICS WASHINGTON, D. C. 20006		2a. REPORT SECURITY CLASSIFICATION
		2b. GROUP
3. REPORT TITLE MULTIWAY CONTINGENCY TABLE ANALYSIS APPLIED TO THE CLASSIFICATION OF MULTIVARIATE DICHOTOMOUS POPULATIONS		
4. DESCRIPTIVE NOTES (Type of report and inclusive dates) TECHNICAL REPORT		
5. AUTHOR(S) (Last name, first name, initial) KULLBACK, S.		
6. REPORT DATE January 9, 1973	7a. TOTAL NO. OF PAGES 22	7b. NO. OF REFS 8
8a. CONTRACT OR GRANT NO. N00014-67-A-0214-0015	8b. ORIGINATOR'S REPORT NUMBER(S) #4	
8c. PROJECT NO. NR-042-267	8d. OTHER REPORT NO(S) (Any other numbers that may be assigned this report)	
10. AVAILABILITY/LIMITATION NOTICES Unlimited. Reproduction in whole or in part is permitted for any purpose of the United States Government.		
11. SUPPLEMENTARY NOTES	12. SPONSORING MILITARY ACTIVITY Office of Naval Research Statistics & Probability Program Arlington, Va. 22217	
13. ABSTRACT <p>Multiway contingency tables, or cross-classifications of vectors of discrete random variables, provide a useful approach to the analysis of multivariate discrete data. In the particular application we shall consider herein, the individual variates are dichotomous or binary. We shall use techniques and concepts presented and discussed by the author in previous papers. We note that the procedures and analysis are not restricted to dichotomous or binary data but are also applicable to polychotomous variates. The procedure we shall use is based on the principle of minimum discrimination information estimation applied to the analysis of multiway contingency tables. It yields results practically equivalent to procedures proposed by other investigators. When the minimum discrimination information estimates provide a satisfactory fit to a set of data, a complete analysis, including significance tests and estimates describing the pattern of observations is provided.</p>		

DD FORM 1473
JAN 64

Unclassified
Security Classification

I

DECLASSIFIED
Security Classification

14. KEY WORDS	LINK A		LINK B		LINK C	
	ROLE	WT	ROLE	WT	ROLE	WT
<p>CONTINGENCY TABLES</p> <p>DICHOTOMOUS POPULATIONS</p> <p>MULTIVARIATE</p>						

INSTRUCTIONS

1. **ORIGINATING ACTIVITY:** Enter the name and address of the contractor, subcontractor, grantee, Department of Defense activity or other organization (corporate author) issuing the report.
- 2a. **REPORT SECURITY CLASSIFICATION:** Enter the overall security classification of the report. Indicate whether "Restricted Data" is included. Marking is to be in accordance with appropriate security regulations.
- 2b. **GROUP:** Automatic downgrading is specified in DoD Directive 5200.10 and Armed Forces Industrial Manual. Enter the group number. Also, when applicable, show that optional markings have been used for Group 3 and Group 4 as authorized.
3. **REPORT TITLE:** Enter the complete report title in all capital letters. Titles in all cases should be unclassified. If a meaningful title cannot be selected without classification, show title classification in all capitals in parentheses immediately following the title.
4. **DESCRIPTIVE NOTES:** If appropriate, enter the type of report, e.g., interim, progress, summary, annual, or final. Give the inclusive dates when a specific reporting period is covered.
5. **AUTHOR(S):** Enter the name(s) of author(s) as shown on or in the report. Enter last name, first name, middle initial. If military, show rank and branch of service. The name of the principal author is an absolute minimum requirement.
6. **REPORT DATE:** Enter the date of the report as day, month, year, or month, year. If more than one date appears on the report, use date of publication.
- 7a. **TOTAL NUMBER OF PAGES:** The total page count should follow normal pagination procedures, i.e., enter the number of pages containing information.
- 7b. **NUMBER OF REFERENCES:** Enter the total number of references cited in the report.
- 8a. **CONTRACT OR GRANT NUMBER:** If appropriate, enter the applicable number of the contract or grant under which the report was written.
- 8b, 8c, & 8d. **PROJECT NUMBER:** Enter the appropriate military department identification, such as project number, subproject number, system number, task number, etc.
- 9a. **ORIGINATOR'S REPORT NUMBER(S):** Enter the official report number by which the document will be identified and controlled by the originating activity. This number must be unique to this report.
- 9b. **OTHER REPORT NUMBER(S):** If the report has been assigned any other report numbers (either by the originator or by the sponsor), also enter this number(s).
10. **AVAILABILITY/LIMITATION NOTICES:** Enter any limitations on further dissemination of the report, other than those

imposed by security classification, using standard statements such as:

- (1) "Qualified requesters may obtain copies of this report from DDC."
- (2) "Foreign announcement and dissemination of this report by DDC is not authorized."
- (3) "U. S. Government agencies may obtain copies of this report directly from DDC. Other qualified DDC users shall request through _____."
- (4) "U. S. military agencies may obtain copies of this report directly from DDC. Other qualified users shall request through _____."
- (5) "All distribution of this report is controlled. Qualified DDC users shall request through _____."

If the report has been furnished to the Office of Technical Services, Department of Commerce, for sale to the public, indicate this fact and enter the price, if known.

11. **SUPPLEMENTARY NOTES:** Use for additional explanatory notes.

12. **SPONSORING MILITARY ACTIVITY:** Enter the name of the departmental project office or laboratory sponsoring (paying for) the research and development. Include address.

13. **ABSTRACT:** Enter an abstract giving a brief and factual summary of the document indicative of the report, even though it may also appear elsewhere in the body of the technical report. If additional space is required, a continuation sheet shall be attached.

It is highly desirable that the abstract of classified reports be unclassified. Each paragraph of the abstract shall end with an indication of the military security classification of the information in the paragraph, represented as (TS), (S), (C), or (U).

There is no limitation on the length of the abstract. However, the suggested length is from 150 to 225 words.

14. **KEY WORDS:** Key words are technically meaningful terms or short phrases that characterize a report and may be used as index entries for cataloging the report. Key words must be selected so that no security classification is required. Identifiers, such as equipment model designation, trade name, military project code name, geographic location, may be used as key words but will be followed by an indication of technical context. The assignment of links, roles, and weights is optional.

II

Multiway Contingency Table Analysis Applied to the Classification of Multivariate Dichotomous Populations

by
S. Kullback

Introduction

Multiway contingency tables, or cross-classifications of vectors of discrete random variables, provide a useful approach to the analysis of multivariate discrete data. In the particular application we shall consider herein, the individual variates are dichotomous or binary. We shall use techniques and concepts presented and discussed in [4] and [6]. We note that the procedures and analysis are not restricted to dichotomous or binary data but are also applicable to polychotomous variates.

For background on the study and problem which gave rise to the data we shall analyze see [8]. In [3], procedures further developed in [4] and [6], were applied to problems of multivariate binary data in information systems, such as communication, pattern recognition, and learning systems. In [1] there is a review of methods and models for the analysis of multivariate binary data. Solomon's data, which we shall analyze herein, is given as a typical example. In [7] there is developed a model based on a set of orthogonal polynomials and applied to Solomon's data. We remark that the procedure we shall use, based on the principle of minimum discrimination information estimation applied to the analysis of multiway contingency tables yields a result practically equivalent to that in [7].

"Multivariate data analysis needs a large and flexible class of hypothetical distributions of free variables indexed by the values of fixed variables. From this class, appropriate subfamilies would be chosen for fitting to specific data sets" [2]. The principle of minimum discrimination information estimation, and its basis, the minimum discrimination information theorem which is quite general in its formulation, lead to exponential families of distributions [4], [5], [6]. The exponential families have very useful and desirable statistical properties and contain many subfamilies in common use [2]. "The data analytic attitude to models is empirical rather than theoretical... when detailed theoretical understanding is unavailable, a more empirical attitude is natural, so that estimation of parameters in models should be seen less as attempts to discover underlying truth and more as data calibrating devices which make it easier to conceive of noisy data in terms of smooth distributions and relations. Exponential families are viewed here as intended for use in the empirical mode. With a given data set, a variety of models may be tried on, and one selected on the ground of looks and fit" [2]. When the minimum discrimination information estimates provide a satisfactory fit to a set of data, a complete analysis, including significance tests and estimates describing the pattern of observations is provided.

Solomon's Data

A total of 2982 high-school seniors were given an attitude questionnaire to assess their attitude towards science. The students were also

classified on the basis of an IQ test into high IQ, the upper half, and low IQ, the lower half. The sixteen possible response vectors to each of four agree-disagree responses were tabulated. The data is given in table 1, where x_1, x_2, x_3, x_4 indicate the statements ([8,p.416]), agree and disagree were coded as 1 and 0 respectively, and listed as low IQ and high IQ. The problem of interest was to determine whether the response vectors could be used as a basis for classifying the students into one of two classes and evaluate possible classification procedures.

Contingency Table Analysis

We shall treat the data as a five-way $2 \times 2 \times 2 \times 2 \times 2$ contingency table, denoting the original observations by $x(hijkl)$, where

$h=1$, low IQ, $h=2$, high IQ ;

$i=1$, response to x_1 coded 0, $i=2$, response to x_1 coded 1;

$j=1$, response to x_2 coded 0, $j=2$, response to x_2 coded 1;

$k=1$, response to x_3 coded 0, $k=2$, response to x_3 coded 1;

$l=1$, response to x_4 coded 0, $l=2$, response to x_4 coded 1.

As a first overview of the data to determine the marginals and their related interaction parameters which may be considered to furnish significant values in the log-linear representation of the exponential family of the estimates [6], we list in table 2a, Analysis of Information, a sequential study of interaction and effect type measures [4], [6].

We remark that the first estimate is

$$x_a^*(hijkl) = x(h \cdots i j k l) / n$$

and the minimum discrimination information statistic (interaction type measure)

$$2I(x:x_a^*) = 2\sum\sum\sum x(hijkl) \ln \frac{x(hijkl)n}{x(h\cdots) x(\cdot i j k l)}$$

tests a null hypothesis that the IQ groupings are homogeneous over the sixteen response vectors [5, Chap.8], [4]. This null hypothesis is rejected and the subsequent study of effect and interaction type measures is an attempt to get a good fit to the data and account for the variation. Although the association between IQ and the response to the first statement is not significant, $2I(x_b^*:x_a^*) = 2.376$, 1 D.F., it was decided to examine in detail the estimate $x_e^*(hijkl)$ whose numerical values are given in table 1. (We remark that it may be shown that

$$2I(x_b^*:x_a^*) = 2\sum\sum x(hi\cdots) \ln \frac{x(hi\cdots)n}{x(h\cdots) x(\cdot i \cdots)}$$

and tests a null hypothesis that IQ is homogeneous over the response to the first question). The estimate $x_e^*(hijkl)$ was selected because it does not differ significantly from the observed values, $2I(x:x_e^*) = 16.307$, 11 D.F. (represents an acceptable fit), is symmetric with respect to the four statements, and is comparable to the first-order model estimate of [7], whose values are also listed in table 1.

From the log-linear representation in figure 1 [6], we obtain the parametric representation for the log-odds (low IQ/high IQ)

$$\ln(x_e^*(1ijkl)/x_e^*(2ijkl))$$

over the sixteen response vectors as given in table 3a. Thus, for example

$$\ln \frac{x_e^*(11111)}{x_e^*(21111)} = \tau_1^h + \tau_{11}^{h1} + \tau_{11}^{hj} + \tau_{11}^{hk} + \tau_{11}^{hl} ,$$

that is, a linear regression of the log-odds in terms of a constant τ_1^h and the main effects of each component of the response vector, namely, $\tau_{11}^{h1}, \tau_{11}^{hj}, \tau_{11}^{hk}, \tau_{11}^{hl}$. The numerical values of the log-odds and the parameters are easily obtained from the entries in the computer output and are also given in table 3a [6].

We note from table 3a that

$$\ln \frac{x_e^*(11jkl)}{x_e^*(21jkl)} - \ln \frac{x_e^*(11jk2)}{x_e^*(21jk2)} = \tau_{11}^{hl} = 0.3338 ,$$

that is, a change from disagree to agree on the fourth statement is associated with an increase of 0.3338 in the log-odds (low IQ/high IQ). Note also that τ_{11}^{hl} represents the association between IQ and response to the fourth statement as measured by the log-cross-product - ratio

$$\tau_{11}^{hl} = \ln \frac{x_e^*(11jkl)x_e^*(21jk2)}{x_e^*(21jkl)x_e^*(11jk2)}$$

and is the same for all eight levels of the responses to statements one, two and three.

Similarly, it is found that

$$\ln \frac{x_e^*(11j1l)}{x_e^*(21j1l)} - \ln \frac{x_e^*(11j2l)}{x_e^*(21j2l)} = \tau_{11}^{hk} = 0.3411 ,$$

$$\ln \frac{x_e^*(11lkl)}{x_e^*(21lkl)} - \ln \frac{x_e^*(112kl)}{x_e^*(212kl)} = \tau_{11}^{hj} = 0.1240 ,$$

$$\ln \frac{x_e^*(11jkl)}{x_e^*(21jkl)} - \ln \frac{x_e^*(12jkl)}{x_e^*(22jkl)} = \tau_{11}^{h1} = -0.2030 .$$

Classification

Since $x(1\cdots) = x_e^*(1\cdots) = 1491$, and $x(2\cdots) = x_e^*(2\cdots) = 1491$, we assign a response vector $(ijk\ell)$ to the region

E_1 : classify as population $h=1$ (low IQ), when

$$\ln \frac{x_e^*(1ijk\ell)}{x_e^*(2ijk\ell)} \geq 0$$

and to the complementary region

E_2 : classify as population $h=2$ (high IQ), when

$$\ln \frac{x_e^*(1ijk\ell)}{x_e^*(2ijk\ell)} < 0 .$$

If we set

$$\mu_1(E_1) = \sum_{(ijk\ell) \in E_1} \frac{x_e^*(1ijk\ell)}{1491} , \quad \mu_2(E_1) = \sum_{(ijk\ell) \in E_1} \frac{x_e^*(2ijk\ell)}{1491} ,$$

then the probability of error of the classification procedure is
[5, pp.4,69,80]

$$\text{Prob Error} = p\mu_2(E_1) + q\mu_1(E_2) = (\mu_2(E_1) + \mu_1(E_2))/2$$

since here $p = x(2\cdots)/2982 = \frac{1}{2}$, $q = x(1\cdots)/2982 = \frac{1}{2}$.

The relevant computations with $x_e^*(hijk\ell)$ are given in table 4(b) and show that the Prob. Error = 0.444. The corresponding computations with the original data $x(hikj\ell)$ are given in table 4(a) and yield Prob. Error = 0.441.

Other Estimates

In view of the measure of the effect of the marginal $x(hi\cdots\ell)$ (and the associated interaction parameters) in table 2a, $2I(x_m^*:x_g^*) = 4.316$, 1D.F.

and the marginal $x(h \cdot j \cdot l)$, $2I(x_p^*: x_n^*) = 3.181$, 1 D.F., the estimate $x_v^*(hijk)$ fitting the marginals $x(\cdot ijk)$, $x(h \cdot j \cdot \cdot)$, $x(h \cdot \cdot k \cdot)$, $x(hi \cdot \cdot l)$ and the estimate $x_w^*(hijk)$ fitting the marginals $x(\cdot ijk)$, $x(h \cdot \cdot k \cdot)$, $x(hi \cdot \cdot l)$, $x(h \cdot j \cdot l)$ were computed. The estimates are given in table 1 and the relevant analysis of information given in table 2b.

The values of the log-odds, parametric representation, and the values of the associated interaction parameters are given in table 3b for $x_v^*(hijk)$ and in table 3c for $x_w^*(hijk)$. Note from table 3b that

$$\ln \frac{x_v^*(11jk1)}{x_v^*(21jk1)} - \ln \frac{x_v^*(11jk2)}{x_v^*(21jk2)} = \tau_{11}^{hl} + \tau_{111}^{hl} = 0.6469 ,$$

$$\ln \frac{x_v^*(12jk1)}{x_v^*(22jk1)} - \ln \frac{x_v^*(12jk2)}{x_v^*(22jk2)} = \tau_{11}^{hl} = 0.2680 ,$$

$$\ln \frac{x_v^*(11jk1)}{x_v^*(21jk1)} - \ln \frac{x_v^*(12jk1)}{x_v^*(22jk1)} = \tau_{11}^{h1} + \tau_{111}^{h1} = -0.0276$$

$$\ln \frac{x_v^*(11jk2)}{x_v^*(21jk2)} - \ln \frac{x_v^*(12jk2)}{x_v^*(22jk2)} = \tau_{11}^{h1} = -0.4065$$

reflecting the interaction of the responses to the first and fourth statements.

From table 3c, it is found for example, that

$$\ln \frac{x_w^*(111k1)}{x_w^*(211k1)} - \ln \frac{x_w^*(111k2)}{x_w^*(211k2)} = \tau_{11}^{hl} + \tau_{111}^{hl} + \tau_{111}^{hjl} = 0.5806$$

$$\ln \frac{x_w^*(121k1)}{x_w^*(221k1)} - \ln \frac{x_w^*(121k2)}{x_w^*(221k2)} = \tau_{11}^{hl} + \tau_{111}^{hjl} = 0.2030$$

$$\ln \frac{x_w^*(112k1)}{x_w^*(212k1)} - \ln \frac{x_w^*(112k2)}{x_w^*(212k2)} = \tau_{11}^{hl} + \tau_{111}^{hl} = 0.9371$$

$$\ln \frac{x_w^*(122k1)}{x_w^*(222k1)} - \ln \frac{x_w^*(122k2)}{x_w^*(222k2)} = \tau_{11}^{hk} = 0.5595$$

reflecting the interactions of the responses to the first, second and fourth statements.

The computation of the probability of error using the estimates $x_v^*(hijk\ell)$ and $x_w^*(hijk\ell)$ is shown in table 4(c) and 4(d) respectively, and yields probabilities of error 0.444 and 0.446.

Measure of Divergence

As a measure of the divergence between the low IQ and high IQ observed and estimated values, we computed the values of

$$J(1,2) = \frac{1}{2} \sum \sum \sum (x(1ijk\ell) - x(2ijk\ell)) \ln \frac{x(1ijk\ell)}{x(2ijk\ell)}$$

for $x(hijk\ell)$, $x_e^*(hijk\ell)$, $x_v^*(hijk\ell)$, $x_w^*(hijk\ell)$ [5, p.130]. The resulting values and their ratios to the respective degrees of freedom are given in table 5. As is to be expected from the properties of the discrimination information we note that

$$J(1,2;x_e^*) < J(1,2;x_v^*) < J(1,2;x_w^*) < J(1,2;x) .$$

However the ratio to the respective degrees of freedom leads to the inequalities

$$J(1,2;x)/D.F. < J(1,2;x_e^*)/D.F. < J(1,2;x_v^*)/D.F. < J(1,2;x_w^*)/D.F.$$

Remark

Martin and Bradley [7] examined Solomon's data in terms of an estimate they called a first-order or linear model. These estimated values are

given in table 1. It turns out that although the underlying approaches are different, the Martin and Bradley parameters and estimates are practically the same as those for $x_e^*(hijkl)$. From [7, pp.216-217] we note that

$$\begin{aligned}\ln \frac{x_e^*(12222)}{x_e^*(22222)} &= \tau_1^h = \ln \frac{1+a_0+a_1+a_2+a_3+a_4}{1-a_0-a_1-a_2-a_3-a_4} \\ \ln \frac{x_e^*(12221)}{x_e^*(22221)} &= \tau_1^h + \tau_{11}^{hl} = \ln \frac{1+a_0+a_1+a_2+a_3-a_4}{1-a_0-a_1-a_2-a_3+a_4} \\ \ln \frac{x_e^*(12212)}{x_e^*(22212)} &= \tau_1^h + \tau_{11}^{hk} = \ln \frac{1+a_0+a_1+a_2-a_3+a_4}{1-a_0-a_1-a_2+a_3-a_4} \\ \ln \frac{x_e^*(12122)}{x_e^*(22122)} &= \tau_1^h + \tau_{11}^{hj} = \ln \frac{1+a_0+a_1-a_2+a_3+a_4}{1-a_0-a_1+a_2-a_3-a_4} \\ \ln \frac{x_e^*(11222)}{x_e^*(21222)} &= \tau_1^h + \tau_{11}^{hi} = \ln \frac{1+a_0-a_1+a_2+a_3+a_4}{1-a_0+a_1-a_2-a_3-a_4}\end{aligned}$$

or to a first approximation

$$\begin{aligned}\tau_1^h &= 2a_0+2a_1+2a_2+2a_3+2a_4 \\ \tau_1^h + \tau_{11}^{hl} &= 2a_0+2a_1+2a_2+2a_3-2a_4 \\ \tau_1^h + \tau_{11}^{hk} &= 2a_0+2a_1+2a_2-2a_3+2a_4 \\ \tau_1^h + \tau_{11}^{hj} &= 2a_0+2a_1-2a_2+2a_3+2a_4 \\ \tau_1^h + \tau_{11}^{hi} &= 2a_0-2a_1+2a_2+2a_3+2a_4.\end{aligned}$$

It is found that

$$\tau_{11}^{hl} = -4a_4$$

$$\tau_{11}^{hk} = -4a_3$$

$$\tau_{11}^{hj} = -4a_2$$

$$\tau_{11}^{hi} = -4a_1 .$$

The values of the parameters given in [7, table 3, p. 217] are

$$a_0 = -0.042, \quad a_1 = 0.049, \quad a_2 = -0.031, \quad a_3 = -0.084, \quad a_4 = -0.082$$

so that

$$\tau_{11}^{hl} = 0.3338 = 0.334, \quad -4a_4 = 0.328$$

$$\tau_{11}^{hk} = 0.3411 = 0.341, \quad -4a_3 = 0.336$$

$$\tau_{11}^{hj} = 0.1240 = 0.124, \quad -4a_2 = 0.124$$

$$\tau_{11}^{hi} = -0.2030 = -0.203, \quad -4a_1 = -0.196$$

The computation for the probability of error using the estimates in [7] are shown in table 4(e) and yields a probability of error 0.445. (Martin and Bradley give a value of the risk as 0.455).

Acknowledgment

This report was prepared under Navy Contract N0014-67-A0214-0015 and partially supported by the Air Force Office of Scientific Research, Office of Aerospace Research. U.S. Air Force under Grant AFOSR-72-2348.

	1 2 3 4 5 6	7 8 9 10 11 12 13 14 15 16	17 18 19 20 21 22 23 24 25 26	27 28 29 30 31	32
h i j k l	h i j k l	h i h j h k h l i j i k i l j k j l k l	h i j h i k h i l h j k h j l h k l i j k i j l i k l j k l	h h h h i	h
1 1 1 1 1	1 1 1 1 1	1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1	1 1 1 1 1	1
1 1 1 1 2	1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1	
1 1 1 2 1	1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1	
1 1 1 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 1 2 1 1	1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1	
1 1 2 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 1 2 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 1 2 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 1 1 1	1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1	
1 2 1 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 1 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 1 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 2 1 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 2 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 2 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
1 2 2 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 1 1 1	1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1	
2 1 1 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 1 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 1 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 2 1 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 2 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 2 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 1 2 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 1 1 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 1 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 1 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 1 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 2 1 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 2 1 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 2 2 1	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
2 2 2 2 2	1 1 1 1 1	1 1 1 1 1	1 1 1 1 1		
x	✓	✓	✓	✓	✓
x ₀	✓	✓	✓	✓	✓
x ₁	✓	✓	✓	✓	✓
x ₂	✓	✓	✓	✓	✓
x ₃	✓	✓	✓	✓	✓

Figure 1

Solomon's Data-Classification Procedures

$x_1 x_2 x_3 x_4$	1j 1k	Observed Low IQ $x(11jk)$	Martin & Bradley	Estimates		Observed High IQ $x(21jk)$	Martin & Bradley	Estimates		
				$x_e^*(11jk)$	$x_v^*(11jk)$			$x_e^*(21jk)$	$x_v^*(21jk)$	$x_v^*(21jk)$
11 11	22 22	62	74.56	74.589	76.097	122	109.45	109.414	107.904	113.844
11 10	22 21	70	67.30	67.296	66.198	68	70.71	70.703	71.802	66.400
11 01	22 12	31	31.32	31.329	31.943	33	32.68	32.671	32.057	34.173
11 00	22 11	41	37.74	37.780	37.337	25	28.26	28.219	28.662	26.115
10 11	21 22	283	266.76	266.570	271.120	329	345.24	345.429	340.879	336.820
10 10	21 21	253	259.17	259.322	254.876	247	240.83	241.675	245.125	249.232
10 01	21 12	200	193.45	193.625	196.841	172	176.55	178.376	175.160	171.963
10 00	21 11	305	314.50	314.491	310.589	217	207.50	207.508	211.411	215.252
01 11	12 22	14	12.10	12.156	10.866	20	21.90	21.844	23.135	24.085
01 10	12 21	11	9.20	9.182	9.929	10	11.80	11.818	11.071	10.240
01 01	12 12	11	9.68	9.659	8.776	11	12.32	12.341	13.224	13.898
01 00	12 11	14	12.02	12.010	12.855	9	10.98	10.990	10.144	9.244
00 11	11 22	31	33.63	33.623	30.125	56	53.37	53.375	56.874	56.179
00 10	11 21	46	47.37	47.263	50.789	55	53.63	53.737	50.211	50.999
00 01	11 12	37	47.54	47.450	43.233	64	53.46	53.550	57.767	56.837
00 00	11 11	82	74.67	74.656	79.426	53	60.33	60.346	55.574	56.517
		1491				1791				

Table 1

Table 2a
Analysis of Information

Marginals Fitted	Information	D.F.
a) $x(.ijk\ell), x(h....)$	$2I(x:x_a^*) = 68.369$	15
b) $x(.ijk\ell), x(hi....)$	$2I(x_b^*:x_a^*) = 2.376$ $2I(x:x_b^*) = 65.993$	1 14
c) $x(.ijk\ell), x(hi....), x(h.j..)$	$2I(x_c^*:x_b^*) = 4.265$ $2I(x:x_c^*) = 61.728$	1 13
d) $x(.ijk\ell), x(hi....), x(h.j..), x(h..k.)$	$2I(x_d^*:x_c^*) = 25.230$ $2I(x:x_d^*) = 36.498$	1 12
e) $x(.ijk\ell), x(hi....), x(h.j..), x(h..k.), x(h...l)$	$2I(x_e^*:x_d^*) = 20.191$ $2I(x:x_e^*) = 16.307$	1 11
f) $x(.ijk\ell), x(h..k.), x(h...l), x(hij..)$	$2I(x_f^*:x_e^*) = 3.016$ $2I(x:x_f^*) = 13.291$	1 10
g) $x(.ijk\ell), x(h...l), x(hij..), x(hi.k.)$	$2I(x_g^*:x_f^*) = 0.042$ $2I(x:x_g^*) = 13.249$	1 9
m) $x(.ijk\ell), x(hij..), x(hi.k.), x(hi...l)$	$2I(x_m^*:x_g^*) = 4.316$ $2I(x:x_m^*) = 8.933$	1 8
n) $x(.ijk\ell), x(hij..), x(hi.k.), x(hi...l), x(h.jk.)$	$2I(x_n^*:x_m^*) = 0.983$ $2I(x:x_n^*) = 7.950$	1 7
p) $x(.ijk\ell), x(hij..), x(hi.k.), x(hi...l), x(h.jk.), x(h.j.l)$	$2I(x_p^*:x_n^*) = 3.181$ $2I(x:x_p^*) = 4.769$	1 6
q) $x(.ijk\ell), x(hij..), x(hi.k.), x(hi...l), x(h.jk.), x(h.j.l),$ $x(h..k\ell)$	$2I(x_q^*:x_p^*) = 0.219$ $2I(x:x_q^*) = 4.550$	1 5
r) $x(.ijk\ell), x(hi...l), x(h.j.l), x(h..k\ell), x(hijk.)$	$2I(x_r^*:x_q^*) = 0.346$ $2I(x:x_r^*) = 4.204$	1 4

Analysis of Information (continued)

Marginals Fitted	Information	D.F.
	$2I(x:x_r^*) = 4.204$	4
s) $x(.ijk\ell), x(h..k\ell), x(hijk.), x(hij.\ell)$	$2I(x_s^*:x_r^*) = 2.303$	1
	$2I(x:x_s^*) = 1.901$	3
t) $x(.ijk\ell), x(hijk.), x(hij.\ell), x(hi.k\ell)$	$2I(x_t^*:x_s^*) = 1.375$	1
	$2I(x:x_t^*) = 0.526$	2
u) $x(.ijk\ell), x(hijk.), x(hij.\ell), x(hi.k\ell), x(h.jk\ell)$	$2I(x_u^*:x_t^*) = 0.361$	1
	$2I(x:x_u^*) = 0.165$	1

Table 2b
Analysis of Information

Marginals Fitted	Information	D.F.
e) $x(.ijk\ell), x(hi...), x(h.j..), x(h..k.), x(h...l)$	$2I(x:x_e^*) = 16.307$	11
v) $x(.ijk\ell), x(h.j..), x(h..k.), x(hi...l)$	$2I(x_v^*:x_e^*) = 3.735$	1
	$2I(x:x_v^*) = 12.572$	10
w) $x(.ijk\ell), x(h..k.), x(hi...l), x(h.j.\ell)$	$2I(x_w^*:x_v^*) = 3.443$	1
	$2I(x:x_w^*) = 9.129$	9

$$\text{Log-odds } \ln \frac{x_e^*(1ijkl)}{x_e^*(2ijkl)}$$

ijkl	Parametric representation					log-odds
1111	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{hl}$	0.2128
1112	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$		-0.1210
1121	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$		$+\tau_{11}^{hl}$	-0.1284
1122	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$			-0.4621
1211	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$	$+\tau_{11}^{hl}$	0.0888
1212	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$		-0.2450
1221	τ_1^h	$+\tau_{11}^{hi}$			$+\tau_{11}^{hl}$	-0.2524
1222	τ_1^h	$+\tau_{11}^{hi}$				-0.5861
2111	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{hl}$	0.4158
2112	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$		0.0820
2121	τ_1^h		$+\tau_{11}^{hj}$		$+\tau_{11}^{hl}$	0.0746
2122	τ_1^h		$+\tau_{11}^{hj}$			-0.2592
2211	τ_1^h			$+\tau_{11}^{hk}$	$+\tau_{11}^{hl}$	0.2918
2212	τ_1^h			$+\tau_{11}^{hk}$		-0.0420
2221	τ_1^h				$+\tau_{11}^{hl}$	-0.0494
2222	τ_1^h					-0.3831

$$\tau_1^h = -0.3831, \tau_{11}^{hi} = -0.2030, \tau_{11}^{hj} = 0.1240$$

$$\tau_{11}^{hk} = 0.3411, \tau_{11}^{hl} = 0.3338$$

Table 3a

$$\text{Log-odds} = \ln \frac{x^*(1ijk\ell)}{x^*(2ijk\ell)}$$

$ijk\ell$	Parametric representation						log-odds
1111	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	0.3571
1112	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$			-0.2898
1121	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$		$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	0.0115
1122	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$				-0.6355
1211	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	0.2366
1212	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$			-0.4101
1221	τ_1^h	$+\tau_{11}^{hi}$			$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	-0.1088
1222	τ_1^h	$+\tau_{11}^{hi}$					-0.7557
2111	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$		0.3847
2112	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$			0.1167
2121	τ_1^h		$+\tau_{11}^{hj}$		$+\tau_{11}^{h\ell}$		0.0390
2122	τ_1^h		$+\tau_{11}^{hj}$				-0.2290
2211	τ_1^h			$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$		0.2644
2212	τ_1^h			$+\tau_{11}^{hk}$			-0.0036
2221	τ_1^h				$+\tau_{11}^{h\ell}$		-0.0813
2222	τ_1^h						-0.3492

$$\tau_1^h = -0.3492, \tau_{11}^{hi} = -0.4065, \tau_{11}^{hj} = 0.1203$$

$$\tau_{11}^{hk} = 0.3457, \tau_{11}^{h\ell} = 0.2680, \tau_{111}^{hi\ell} = 0.3789$$

Table 3b

$$\text{Log-odds} = \ln \frac{x_w^*(1ijk\ell)}{x_w^*(2ijk\ell)}$$

$ijk\ell$	Parametric representation							log-odds
1111	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	$+\tau_{111}^{hj\ell}$	0.3283
1112	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$				-0.2523
1121	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$		$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$	$+\tau_{111}^{hj\ell}$	-0.0197
1122	τ_1^h	$+\tau_{11}^{hi}$	$+\tau_{11}^{hj}$					-0.6004
1211	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$		0.3976
1212	τ_1^h	$+\tau_{11}^{hi}$		$+\tau_{11}^{hk}$				0.5396
1221	τ_1^h	$+\tau_{11}^{hi}$			$+\tau_{11}^{h\ell}$	$+\tau_{111}^{hi\ell}$		0.0495
1222	τ_1^h	$+\tau_{11}^{hi}$						-0.8876
2111	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$		$+\tau_{111}^{hj\ell}$	0.3542
2112	τ_1^h		$+\tau_{11}^{hj}$	$+\tau_{11}^{hk}$				0.1512
2121	τ_1^h		$+\tau_{11}^{hj}$		$+\tau_{11}^{h\ell}$		$+\tau_{111}^{hj\ell}$	0.0061
2122	τ_1^h		$+\tau_{11}^{hj}$					-0.1968
2211	τ_1^h			$+\tau_{11}^{hk}$	$+\tau_{11}^{h\ell}$			0.4235
2212	τ_1^h			$+\tau_{11}^{hk}$				-0.1360
2221	τ_1^h				$+\tau_{11}^{h\ell}$			0.0754
2222	τ_1^h							-0.4841

$$\tau_1^h = -0.4841, \tau_{11}^{hi} = -0.4035, \tau_{11}^{hj} = 0.2873$$

$$\tau_{11}^{hk} = 0.3481, \tau_{11}^{h\ell} = 0.5595, \tau_{111}^{hi\ell} = 0.3776$$

$$\tau_{111}^{hj\ell} = -0.3565$$

Table 3c

$E_1: \{ijk: \text{in odds} \geq 0\}$

E_1 : Observations

$E_1: x_e^*$

ijk	$x(i,j,k)$	$x(2ijk)$	ijk	$x_e^*(1ijk)$	$x_e^*(2ijk)$
1111	82	53	1111	74.656	60.346
1211	14	9	1211	12.010	10.990
1221	11	10	2111	314.491	207.508
2111	305	217	2112	193.625	178.376
2112	200	172	2121	259.322	240.679
2121	253	247	2211	$\frac{37.780}{891.884}$	$\frac{28.219}{726.113}$
2211	41	25			
2221	$\frac{70}{976}$	$\frac{68}{801}$			

$$\mu_2(E_1) = \frac{801}{1491}, \quad \mu_1(E_2) = \frac{1491-976}{1491}$$

$$\mu_2(E_1) = \frac{726.118}{1491}, \quad \mu_1(E_2) = \frac{1491-891.884}{1491}$$

$$\text{Prob. Error} = \frac{1}{2} \frac{801+515}{1491}$$

$$\text{Prob. Error} = \frac{1}{2} \frac{726.118+599.116}{1491}$$

$$= \frac{1316}{2 \times 1491} = 0.441$$

$$= \frac{1325.234}{2982}$$

$$= 0.444$$

(a)

(b)

Table 4

$E_1: x_v^*$	$x_v^*(11jkl)$	$x_v^*(21jkl)$
1jkl		
1111	79.426	55.574
1121	50.789	50.211
1211	12.855	10.144
2111	310.589	211.411
2112	196.841	175.160
2121	254.876	245.125
2211	<u>37.327</u>	<u>28.662</u>
	942.713	776.287

$$\mu_2(E_1) = \frac{776.287}{1491}$$

$$\mu_1(E_2) = \frac{1491-942.713}{1491}$$

$$\text{Prob. Error} = \frac{1}{2} \frac{776.287+548.287}{1491}$$

$$= \frac{1324.574}{2982}$$

$$= 0.444$$

(c)

Table 4

$E_1: x_v^*$	$x_v^*(11jkl)$	$x_v^*(21jkl)$
1jkl		
1111	78.482	56.517
1211	13.756	9.244
1212	8.102	13.898
1221	10.760	10.240
2111	306.748	215.252
2112	200.037	171.963
2121	250.769	249.232
2211	39.884	26.115
2221	<u>71.600</u>	<u>66.401</u>
	980.138	818.862

$$\mu_2(E_1) = \frac{818.862}{1491}$$

$$\mu_1(E_2) = \frac{1491-980.138}{1491}$$

$$\text{Prob. Error} = \frac{1}{2} \frac{818.862+510.862}{1491}$$

$$= \frac{1329.724}{2982}$$

$$= 0.446$$

(d)

E_1	$\hat{x}(1ijk\ell)$	$\hat{x}(2ijk\ell)$
1111	74.67	60.33
1211	12.02	10.98
2111	314.50	207.50
2112	193.45	178.55
2121	259.17	240.83
2211	$\frac{37.74}{891.55}$	$\frac{28.26}{726.45}$

$$\mu_2(E_1) = \frac{726.45}{1491}, \quad \mu_1(E_2) = \frac{1491-891.55}{1491}$$

$$\text{Prob. Error} = \frac{1}{2} \frac{726.45+599.45}{1491}$$

$$= \frac{1325.90}{2982}$$

$$= 0.445$$

Table 4(e)

Divergence Between Low IQ and High IQ
Observations and Estimates

$$\frac{1}{2} \sum \sum \sum \sum (x(11jkl) - x(21jkl)) \ln \frac{x(11jkl)}{x(21jkl)} = 69.132$$

$$69.132/15 = 4.61/\text{D.F.}$$

$$\frac{1}{2} \sum \sum \sum \sum (x_{\bullet}^*(11jkl) - x_{\bullet}^*(21jkl)) \ln \frac{x_{\bullet}^*(11jkl)}{x_{\bullet}^*(21jkl)} = 52.374$$

$$52.374/11 = 4.76/\text{D.F.}$$

$$\frac{1}{2} \sum \sum \sum \sum (x_V^*(11jkl) - x_V^*(21jkl)) \ln \frac{x_V^*(11jkl)}{x_V^*(21jkl)} = 56.249$$

$$56.249/10 = 5.62/\text{D.F.}$$

$$\frac{1}{2} \sum \sum \sum \sum (x_W^*(11jkl) - x_W^*(21jkl)) \ln \frac{x_W^*(11jkl)}{x_W^*(21jkl)} = 59.815$$

$$59.815/9 = 6.65/\text{D.F.}$$

Table 5

References

- [1] Cox, D.R. (1972), The analysis of multivariate binary data, Applied Statistics, 21, 113-120.
- [2] Dempster, A.P. (1971), An overview of multivariate data analysis, Journal of Multivariate Analysis, 1, 316-346.
- [3] Ku, H.H. and Kullback, S. (1969), Approximating discrete probability distributions, IEEE Trans. on Information Theory, IT-15, 444-447.
- [4] Ku, H.H., Varner, R.N., and Kullback, S. (1971), On the analysis of multidimensional contingency tables, Journal of the American Statistical Association, 66, 55-64.
- [5] Kullback, S. (1959), Information Theory and Statistics, Wiley, N.Y. 1968 Edition, Dover Publications Inc. N.Y.
- [6] Kullback, S. (1970), Minimum discrimination information estimation and application, Invited paper presented to Sixteenth Conference on the Design of Experiments in Army Research, Development and Testing, 21 October 1970. ARO-D Report 71-3, 1-38 Proceedings of the Conference.
- [7] Martin, D.C. and Bradley, R.A. (1972), Probability models, estimation, and classification for multivariate dichotomous populations, Biometrics, 28, 203-221.
- [8] Solomon, H. (1960), Classification procedures based on dichotomous response vectors, No. 36 in Contributions to Probability and Statistics, Essays in Honor of Harold Hotelling, Edited by I. Olkin et al Stanford U.P., Stanford, Cal. 1960, pp. 414-423 (Also in Studies in Item Analysis and Prediction, Edited by H. Solomon, Stanford U.P. Stanford, Cal, 1961, pp. 177-186.